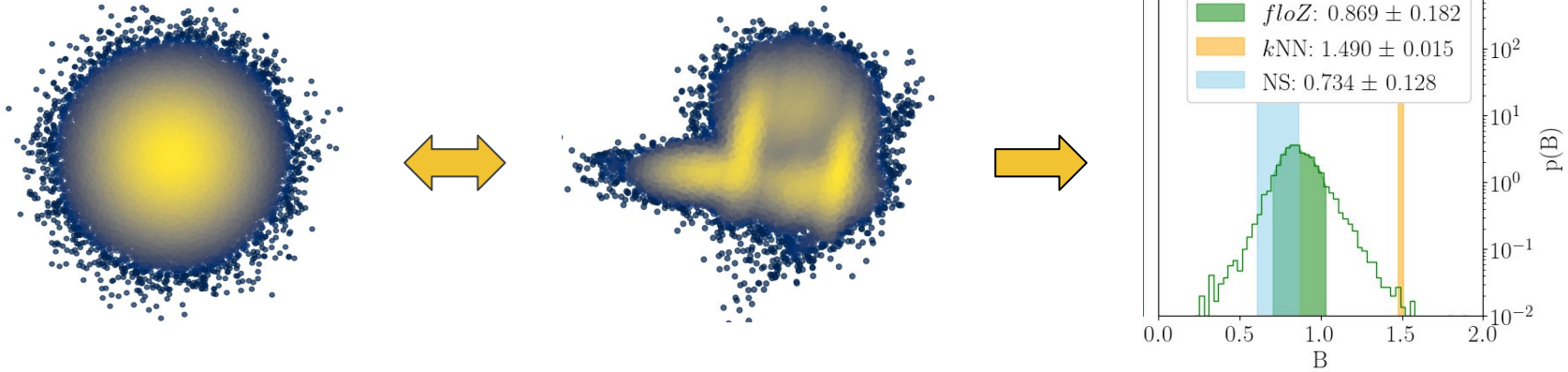


## Bayes factor from normalizing flows

Rahul **Srinivasan**, Marco Crisostomi, Roberto Trotta, Enrico Barausse, and Matteo Breschi



# The Bayes Theorem

$$p(A|B) = \frac{p(B|A) \cdot p(A)}{p(B)}$$

# The Bayes Theorem

$$p(A|B, C) = \frac{p(B|A, C) \cdot p(A|C)}{p(B|C)}$$

# The Bayes Theorem

$$p(\theta|\text{data}, M) = \frac{p(\text{data}|\theta, M) \cdot p(\theta|M)}{p(\text{data}|M)}$$

# The Bayes Theorem

$$\textit{Posterior} = \frac{\textit{Likelihood} \cdot \textit{Prior}}{\textit{Evidence}}$$

# Competing models

$$p(\theta|\text{data}, M_1) = \frac{p(\text{data}|\theta, M_1) \cdot p(\theta|M_1)}{p(\text{data}|M_1)}$$

*Does the data favour  $M_1$  or  $M_2$ ?  
And by how much?*

$$p(\theta|\text{data}, M_2) = \frac{p(\text{data}|\theta, M_2) \cdot p(\theta|M_2)}{p(\text{data}|M_2)}$$

# Competing models

$$p(\theta|\text{data}, M_1) = \frac{p(\text{data}|\theta, M_1) \cdot p(\theta|M_1)}{p(\text{data}|M_1)}$$

$$p(\theta|\text{data}, M_2) = \frac{p(\text{data}|\theta, M_2) \cdot p(\theta|M_2)}{p(\text{data}|M_2)}$$

## The Bayes Factor

*By what factor does the data favour  $M_1$  over  $M_2$ ?*

$$B = \frac{p(\text{data}|M_1)}{p(\text{data}|M_2)}$$
$$= \frac{\text{Evidence}_1}{\text{Evidence}_2}$$

# The Evidence

$$p(\theta|\text{data},M) = \frac{p(\text{data}|\theta,M) \cdot p(\theta|M)}{p(\text{data}|M)}$$

Probability density  
i.e., normalized.



# The Evidence

$$p(\theta|\text{data},M) = \frac{p(\text{data}|\theta,M) \cdot p(\theta|M)}{p(\text{data}|M)}$$

$$p(\text{data}|M) = \int p(\text{data}|\theta,M) \cdot p(\theta|M) d\theta$$

$$\textit{Evidence} = \int \textit{Likelihood} \cdot \textit{Prior} d\theta$$

Computing this integral can be quite non-trivial, and often, intractable.

# State of the art and their *shortcomings*

# State of the art and their *shortcomings*

## Nested sampling<sup>1</sup>:

Evidence estimated by iteratively computing the likelihood.

- ***Computationally intensive*** likelihood *recalculation*.
- ***Slow***, CPU calculations, not parallelizable with GPUs.
- ***Scalability*** issues for high dimensions
  - Ex: 150 dimensions are computationally prohibitive

1. John Skilling “Nested Sampling,” 10.1063/1.1835238.

# State of the art and their *shortcomings*

## Nested sampling<sup>1</sup>:

Evidence estimated by iteratively computing the likelihood.

- *Computationally intensive* likelihood recalculation.
- *Slow*, CPU calculations, not parallelizable with GPUs.
- *Scalability* issues for high dimensions

## Other techniques:

1. k-nearest neighbours<sup>2</sup>, Laplace approx. - *Less expressive*: fails for large non-gaussianity.
2. Normalizing flow-based nested<sup>3</sup>/Gaussianized bridge<sup>4</sup> sampling - *Requires likelihood re-calculation*

1. John Skilling “Nested Sampling,” 10.1063/1.1835238.

2. A. Heavens, et al 2017 arXiv:1704.03472 [stat.CO]

3. Nested sampling with normalizing flows for gravitational-wave inference, 10.1103/PhysRevD.103.103006

4. Jia, He; Seljak, Uroš, 2019 10.48550/arXiv.1912.06073

# State of the art and their *shortcomings*

## Nested sampling<sup>1</sup>:

Evidence estimated by iteratively computing the likelihood.

Likelihood evaluation can be expensive.

These are pre-computed for MCMC samples in parameter estimation pipelines.

Why not use it?

Useful to have a fast, scalable, and expressive method that does not require extra likelihood evaluations.

1. k-nearest neighbours<sup>2</sup>, Laplace approx. - *Less expressive*: fails for large non-gaussianity.
2. Normalizing flow-based nested<sup>3</sup>/Gaussianized bridge<sup>4</sup> sampling - *Requires likelihood re-calculation*

1. John Skilling “Nested Sampling,” 10.1063/1.1835238.

2. A. Heavens, et al 2017 arXiv:1704.03472 [stat.CO]

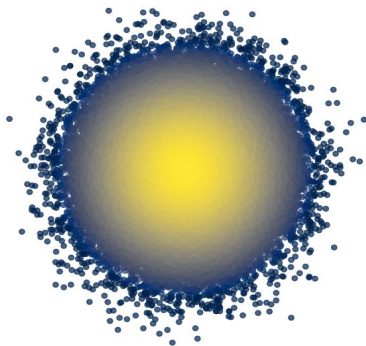
3. Nested sampling with normalizing flows for gravitational-wave inference, 10.1103/PhysRevD.103.103006

4. Jia, He; Seljak, Uroš, 2019 10.48550/arXiv.1912.06073

# A normalizing flow

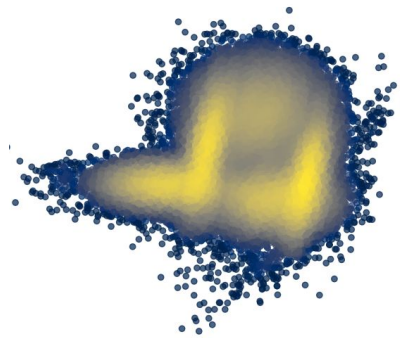
Flows solves for a bijective map b/ the *latent* Normal distribution and the *real* non-trivial distribution.

Known *latent* distribution



$$\mathbf{y} \sim n(\mathbf{y})$$

Target *real* distribution

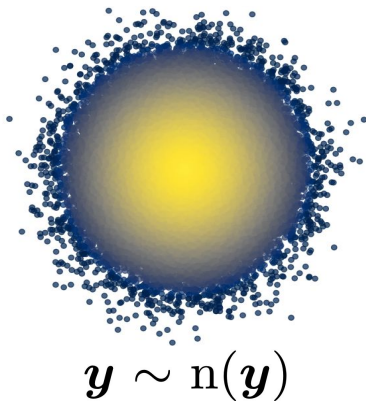


$$\mathbf{x} \sim p(\mathbf{x})$$

# A normalizing flow

Flows solves for a bijective map b/ the *latent* Normal distribution and the *real* non-trivial distribution.

Known *latent* distribution

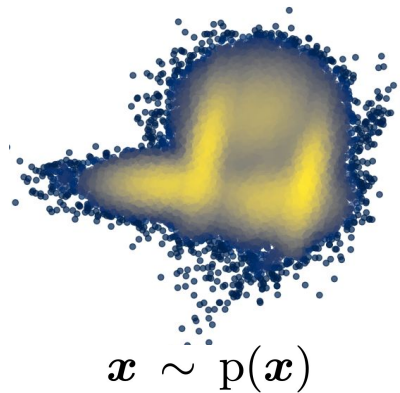


$$\mathbf{x} = \mathbf{f}_\phi(\mathbf{y})$$



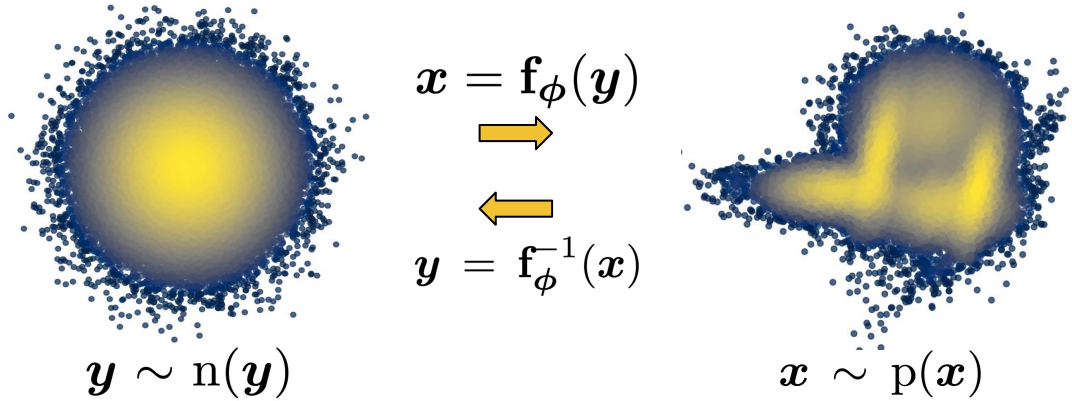
$$\mathbf{y} = \mathbf{f}_\phi^{-1}(\mathbf{x})$$

Target *real* distribution



# A normalizing flow

Flows solves for a bijective map b/ the *latent* Normal distribution and the *real* non-trivial distribution.



$$\text{Target distribution } p(\mathbf{x}) \mapsto q_\phi(\mathbf{x}) \quad \text{Flow prediction} = n(\mathbf{f}_\phi^{-1}(\mathbf{x})) \left| \det \frac{\partial \mathbf{f}_\phi^{-1}}{\partial \mathbf{x}}(\mathbf{x}) \right|$$

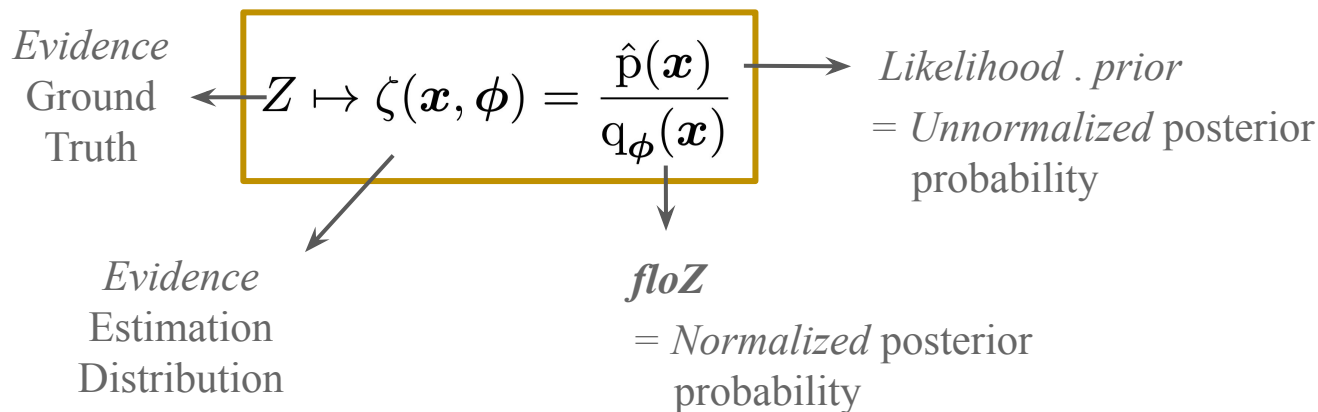


# Theory behind *floZ*

Evidence = normalization constant of likelihood x prior

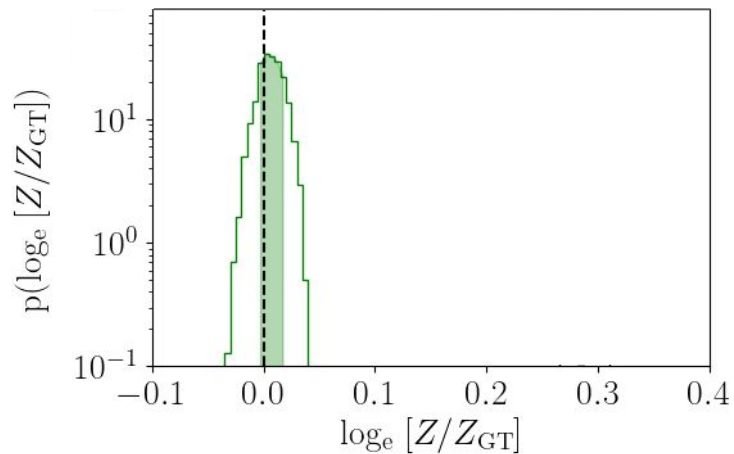
# Theory behind *floZ*

Evidence = normalization constant of likelihood x prior



# Expected output:

Evidence *distribution*



Ideally a delta function

# Implementation:

## Loss Terms

1. Normalizing flow loss:

$$\mathcal{L}_1(\phi) = -\underset{\substack{\downarrow \\ \text{Expectation over posterior samples}}}{\mathbb{E}_{p(\mathbf{x})}} [\log(\underset{\substack{\nearrow \\ \text{floZ prediction}}}{q_\phi(\mathbf{x}))})]$$

# Implementation:

## Loss Terms

1. Normalizing flow loss:

$$\mathcal{L}_1(\phi) = -\underset{\substack{\downarrow \\ \text{Expectation over posterior samples}}}{\mathbb{E}_{p(\mathbf{x})}} [\log(q_\phi(\mathbf{x}))]$$

*↑ floZ prediction*

2. Reducing evidence estimation error:

$$\mathcal{L}_2(\phi) \simeq \log \sigma_{\mathfrak{h}}$$

*↓ Standard deviation of evidence estimation*

# Implementation:

## Loss Terms

1. Normalizing flow loss:

$$\mathcal{L}_1(\phi) = -\underset{\substack{\uparrow \\ \text{Expectation over posterior samples}}}{\mathbb{E}_{p(\mathbf{x})}} [\log(q_\phi(\mathbf{x}))]$$

*floZ prediction*

2. Reducing evidence estimation error:

$$\mathcal{L}_2(\phi) \simeq \log \sigma_{\mathfrak{h}}$$

*Standard deviation  
of evidence estimation*

3. Identity evidence ratio of all pairs of samples:

$$\mathcal{L}_{3a}(\phi) = |\log \mu_{\mathfrak{g}}|$$

*Mean evidence ratio*

# Implementation:

## Loss Terms

1. Normalizing flow loss:

$$\mathcal{L}_1(\phi) = -\underset{\substack{\uparrow \\ \text{floZ prediction}}}{\text{Expectation over posterior samples}} \mathbb{E}_{p(\mathbf{x})} [\log(q_\phi(\mathbf{x}))]$$

2. Reducing evidence estimation error:

$$\mathcal{L}_2(\phi) \simeq \log \sigma_h \underset{\substack{\downarrow \\ \text{Standard deviation} \\ \text{of evidence estimation}}}{}$$

3. Identity evidence ratio of all pairs of samples:

$$\mathcal{L}_{3a}(\phi) = |\log \mu_g| \underset{\substack{\uparrow \\ \text{Mean evidence ratio}}}{}$$

4. Reducing evidence ratio error:

$$\mathcal{L}_{3b}(\phi) = \log \sigma_g \underset{\substack{\downarrow \\ \text{Standard deviation of the} \\ \text{ratio of evidence}}}{}$$

# Implementation:

## Loss Terms

1. Normalizing flow loss:

$$\mathcal{L}_1(\phi) = \text{Expect} \left[ \mathbf{L}_1 \left[ \underset{\text{rior samples}}{\mathcal{Q}_\phi(\mathbf{x})} \right] \right] \overset{\text{floZ prediction}}{\uparrow}$$

2. Reducing evidence estimation error:

$$\mathcal{L}_2(\phi) \simeq 1 \quad \mathbf{L}_2 \quad \text{viation of evidence estimation}$$

3. Identity evidence ratio of all pairs of samples:

$$\mathcal{L}_{3a}(\phi) : \mathbf{L}_{3a} \quad \text{idence ratio}$$

4. Reducing evidence ratio error:

$$\mathcal{L}_{3b}(\phi) = \mathbf{L}_{3b} \quad \text{Standard deviation of the ratio of evidence}$$



# Implementation:

## Loss Scheduling

Solving the four losses simultaneously:

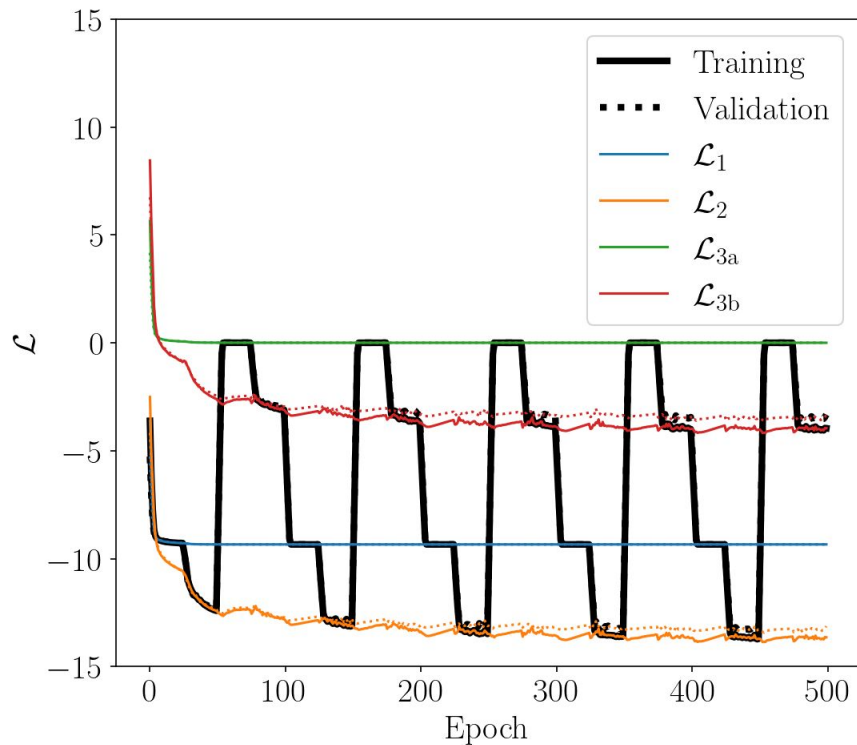
- 1) Weighted sum of losses.
- 2) Schedule the losses

# Implementation:

## Loss Scheduling

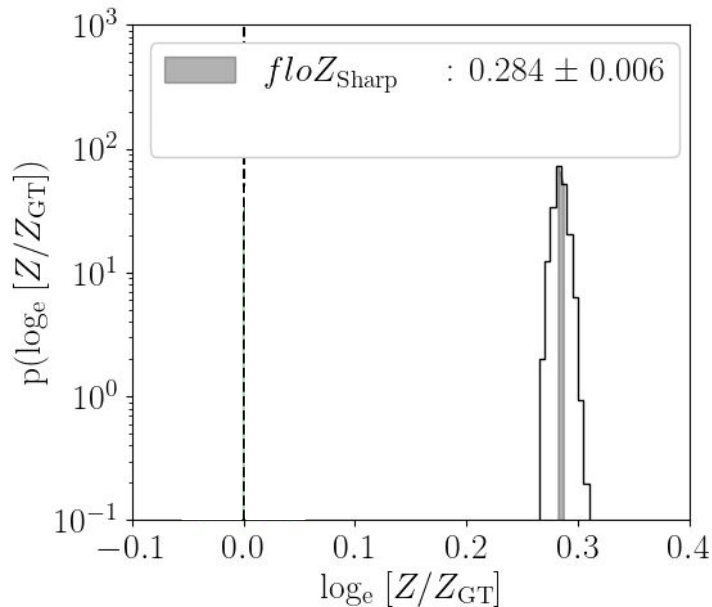
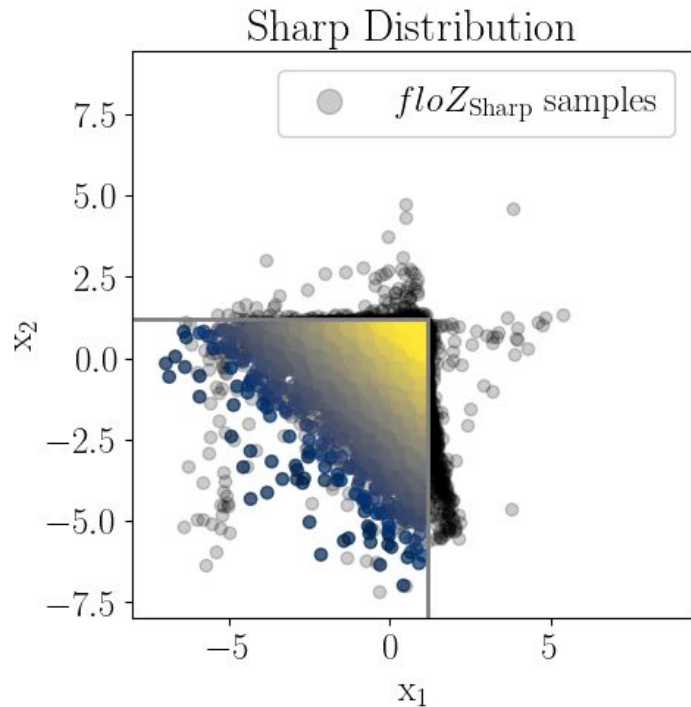
Solving the four losses simultaneously:

- 1) Weighted sum of losses.
- 2) *Schedule the losses*



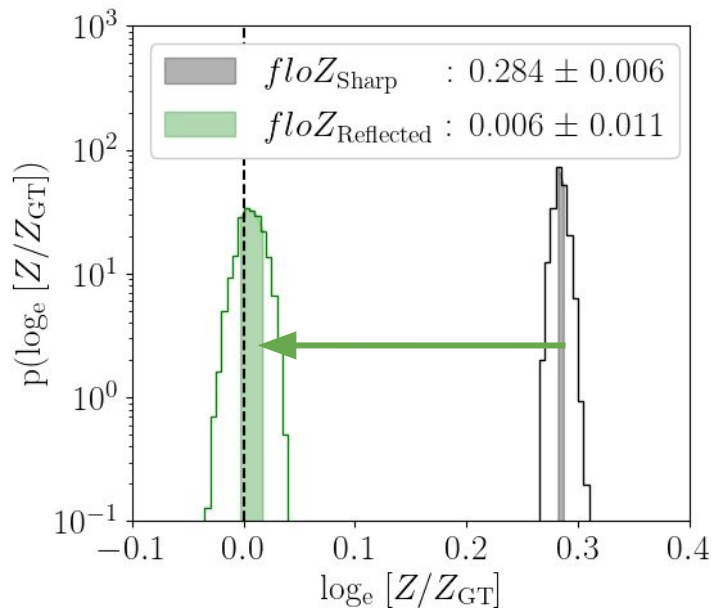
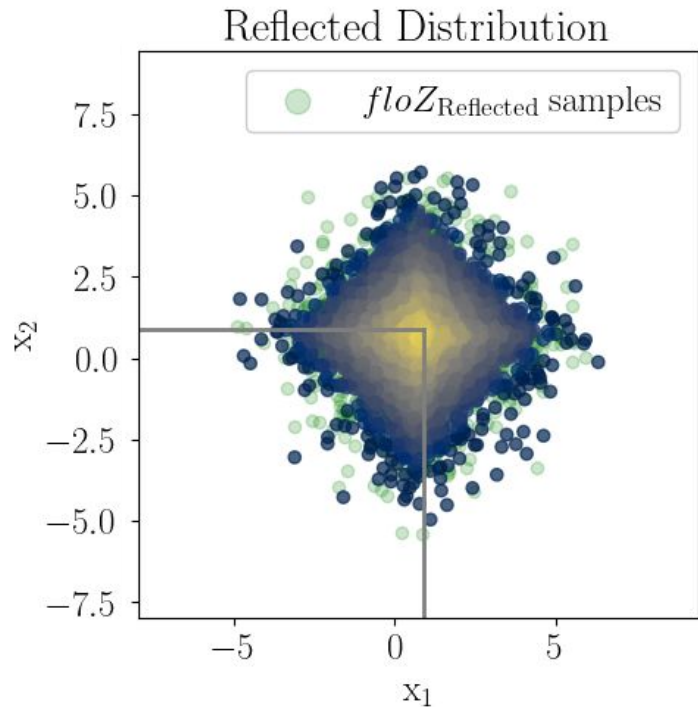
# Implementation:

## Dealing with sharp boundaries



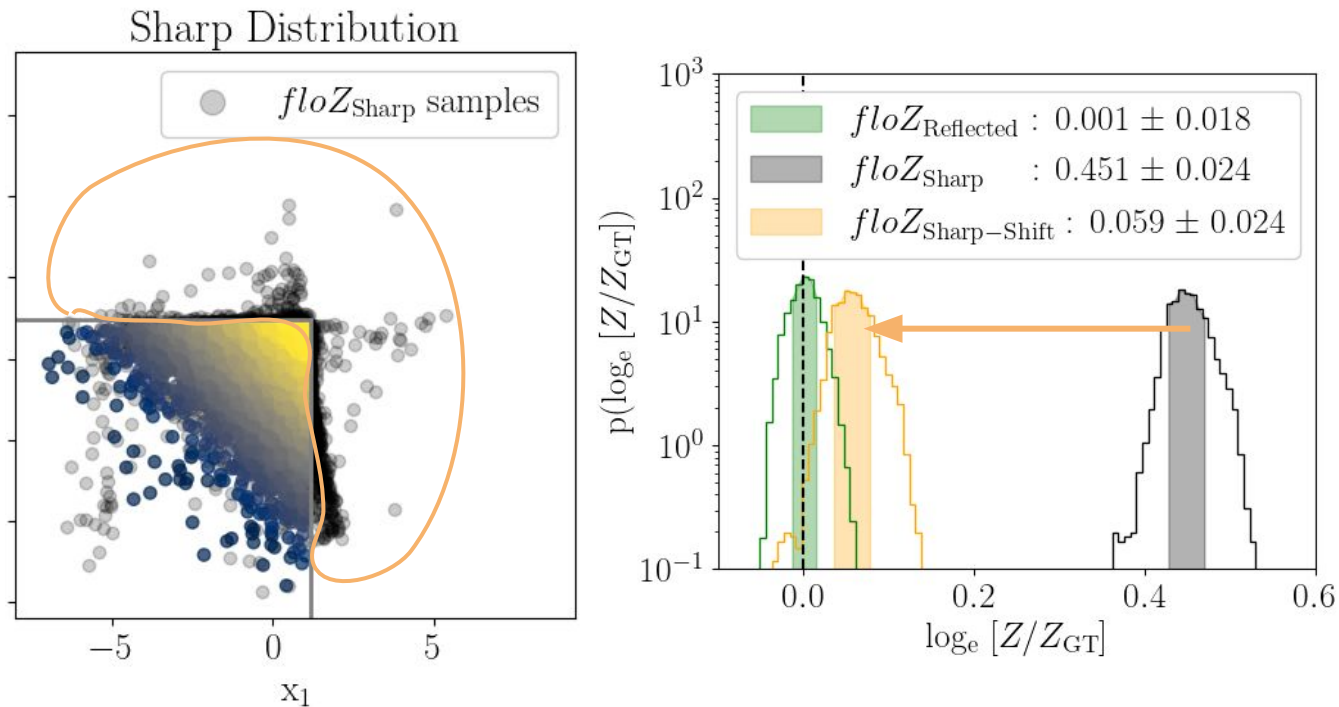
# Implementation:

## Dealing with sharp boundaries



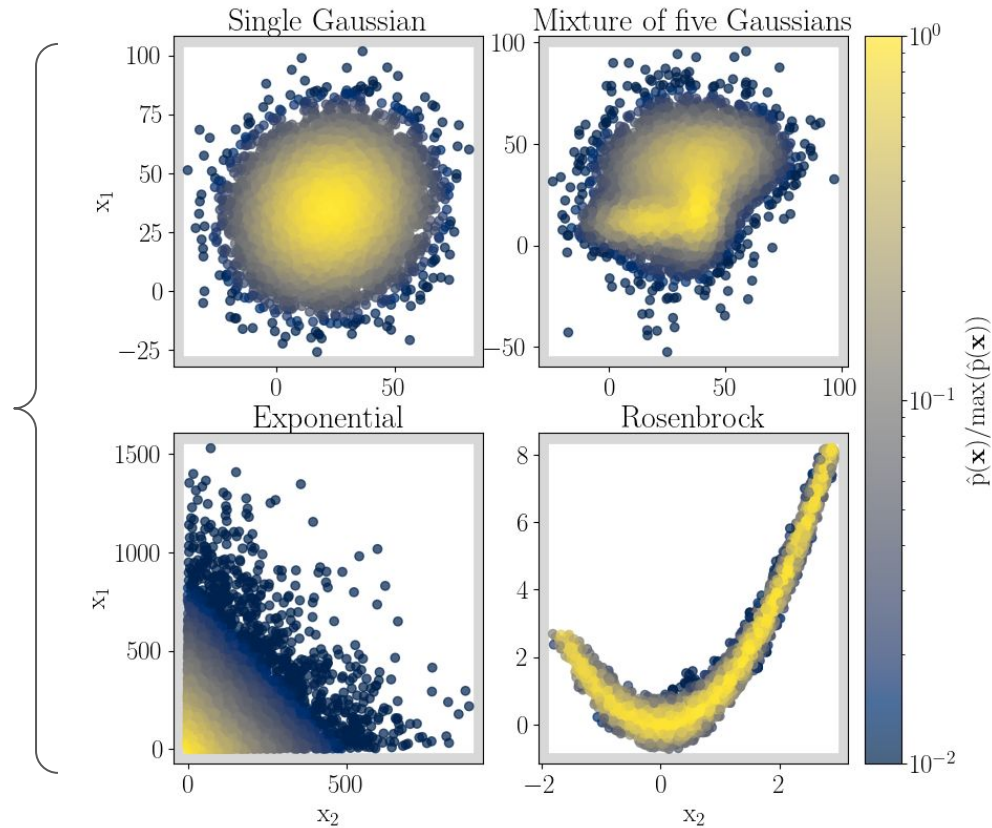
# Alternatives?

Reweighting by fraction of outliers



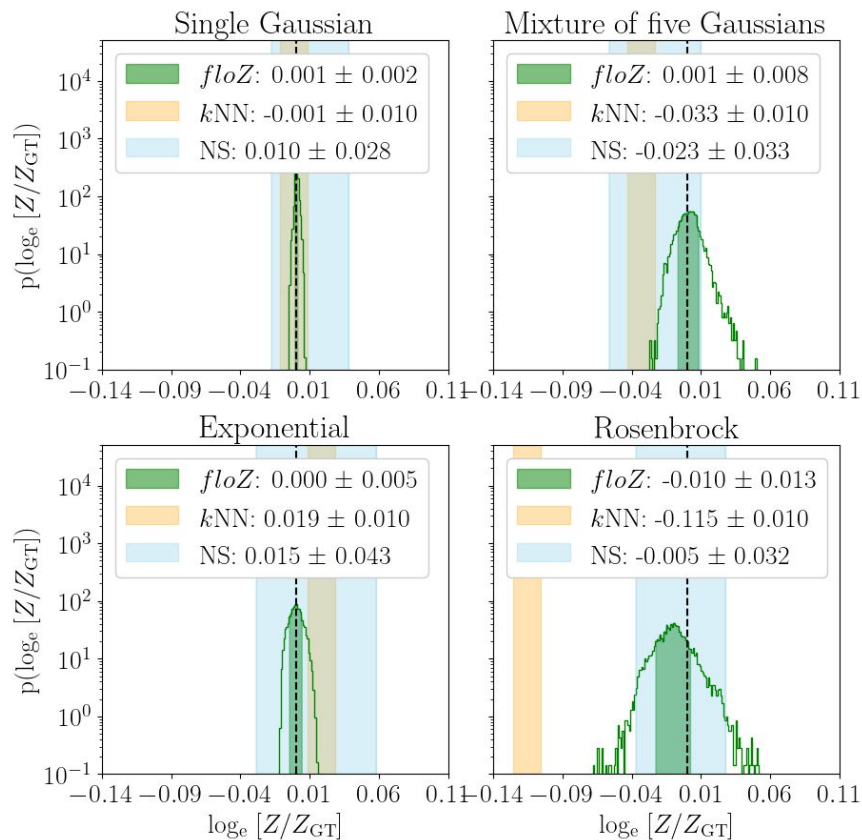
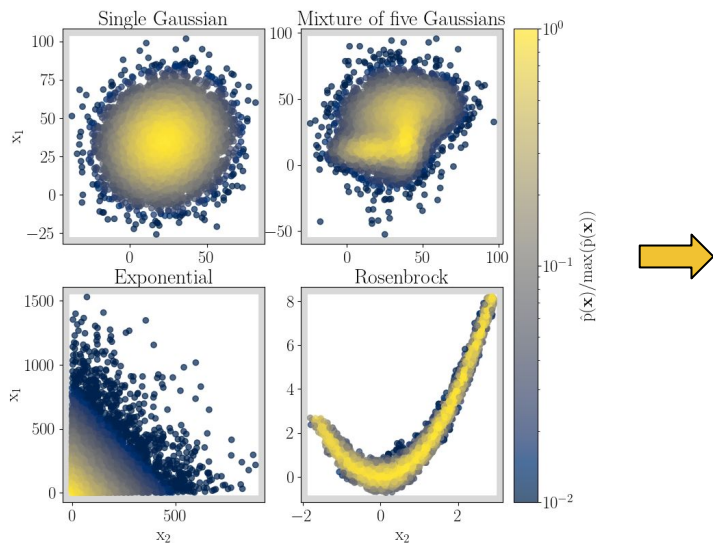
# Distributions for benchmarking

2, 10, 15 dimensions



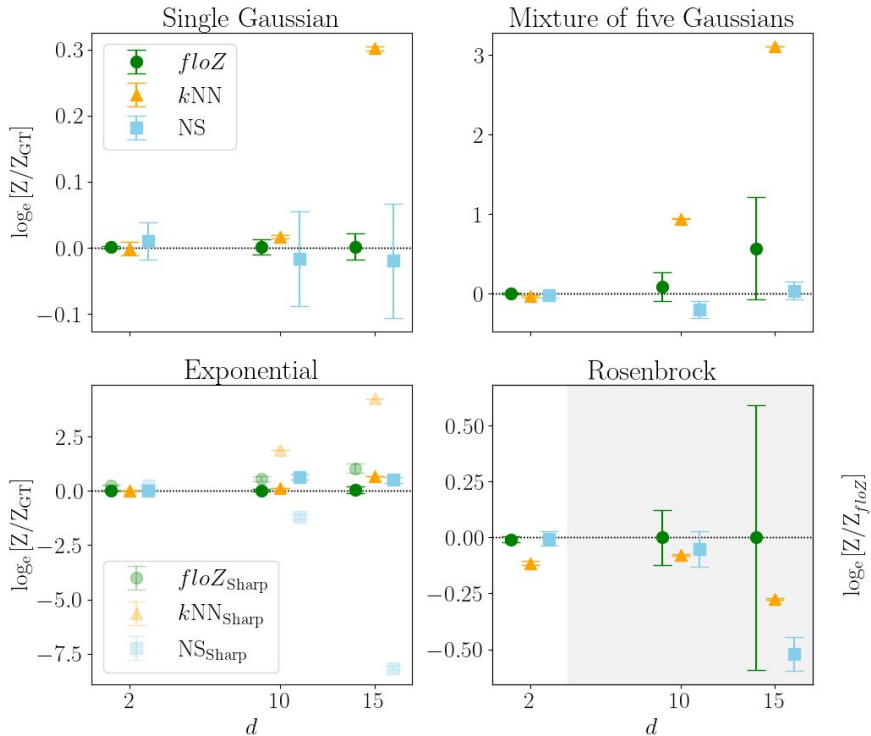
# Benchmarking w/ StateOfTheArt

kNN: k-Nearest Neighbours  
NS: Nested Sampling



# Benchmarking w/ StateOfTheArt

4 Distributions x {2,10,15} Dimensions



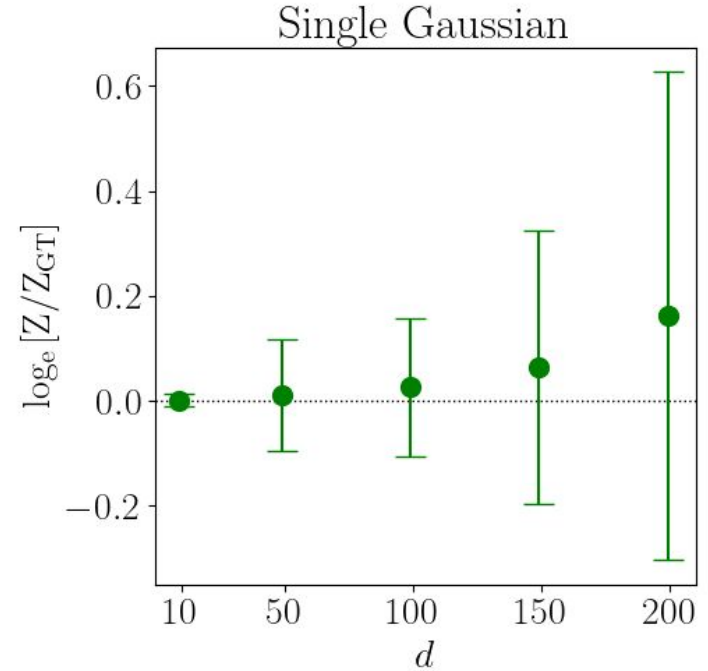
- **Accurate:**  
*floZ* and NS are in good agreement.  
Outperforms *kNN*
- **Scalable:**  
15d require no more than  $10^5$  samples.
- **Rapid**  
15d results of *floZ* obtained in  $\sim 20$ min on an A100 GPU



# High dimensional scalability

For the same number of samples ( $10^5$ ) & model complexity.

\* For complex distributions, we need a combination of more samples, longer training time, and deeper networks.



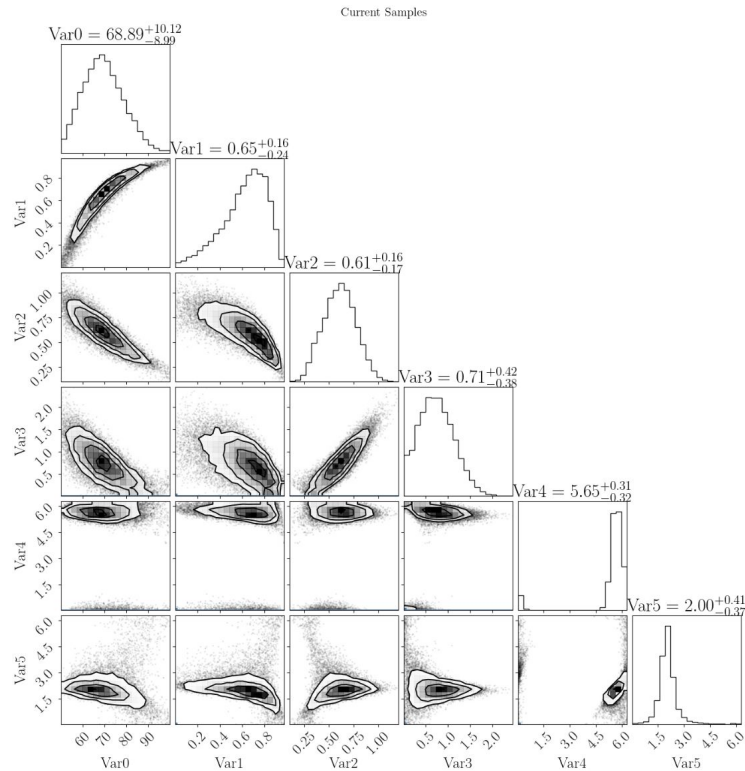
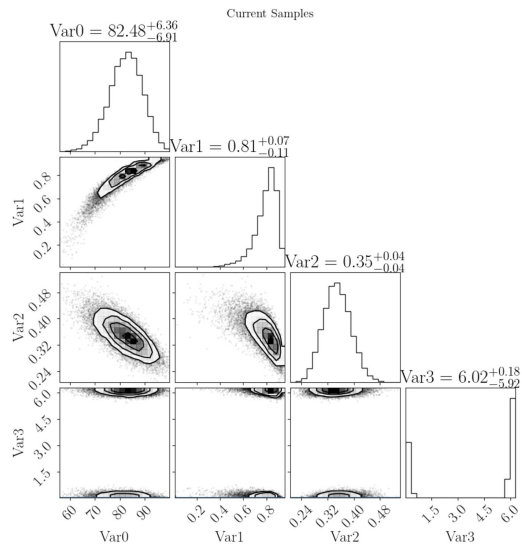
# Applications: GW Ringdown

Bayes factor in favor of the presence of the higher 221 overtone in GW150914

Fundamental Mode

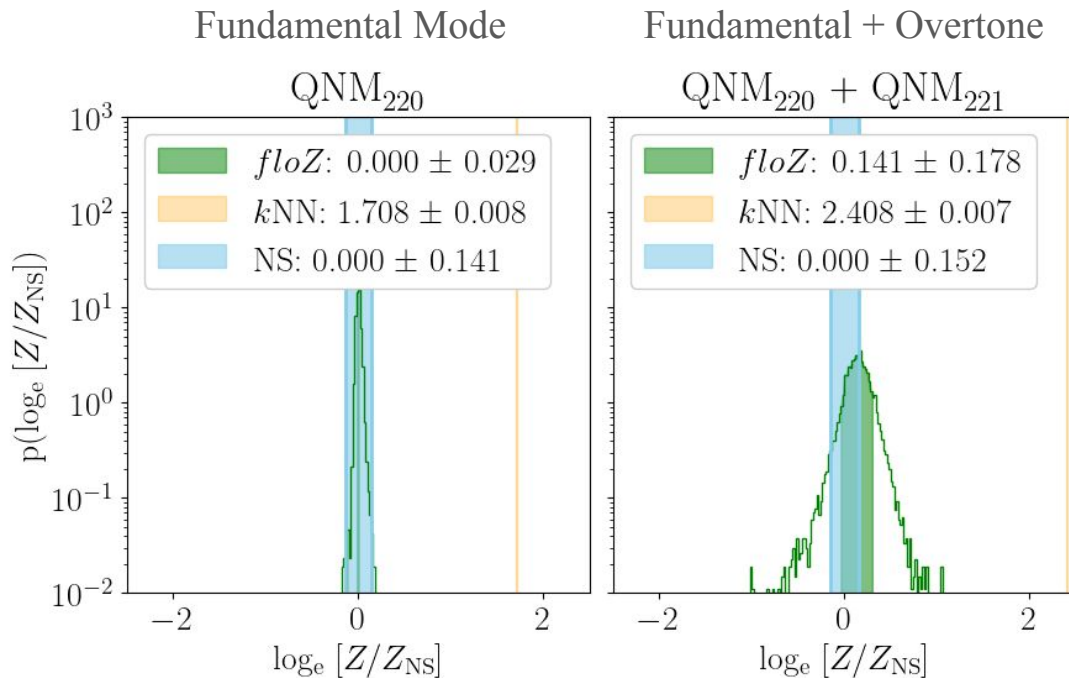
vs

Fundamental Mode w/ Overtone



# Applications: GW Ringdown

Bayes factor in favor of the presence of the higher 221 overtone in GW150914

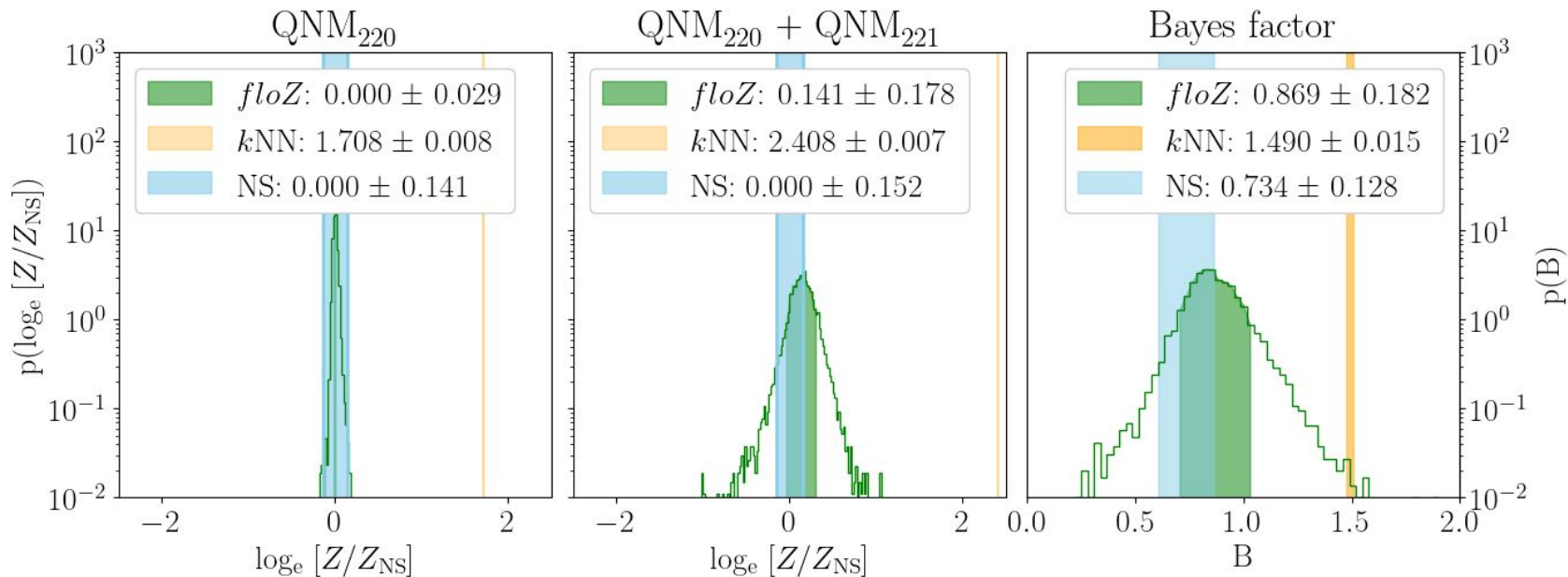


Samples from a nested sampler:  
CPNest<sup>1</sup>

1. W. Del Pozzo and J. Veitch, “CPNest: Parallel nested sampling.” Astrophysics source code library, record ascl:2205.021, May, 2022.

# Applications: GW Ringdown

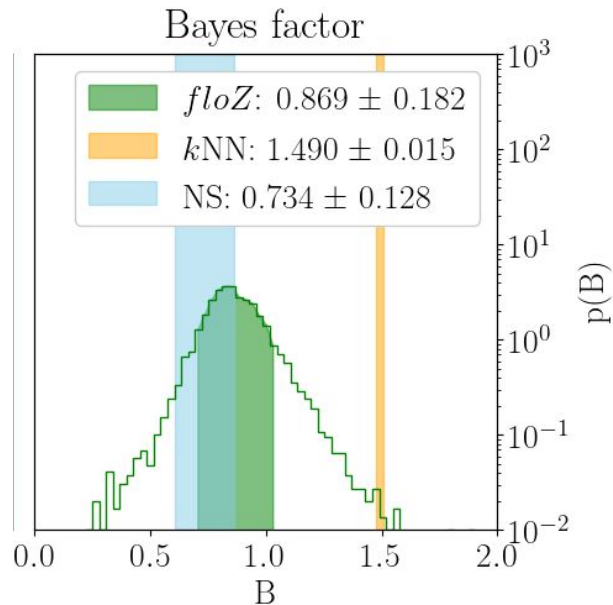
Bayes factor in favor of the presence of the higher 221 overtone in GW150914



# Applications: GW Ringdown

Bayes factor in favor of the presence of the higher 221 overtone in GW150914

floZ estimates is compatible with nested sampling within their  $1\sigma$  uncertainties.



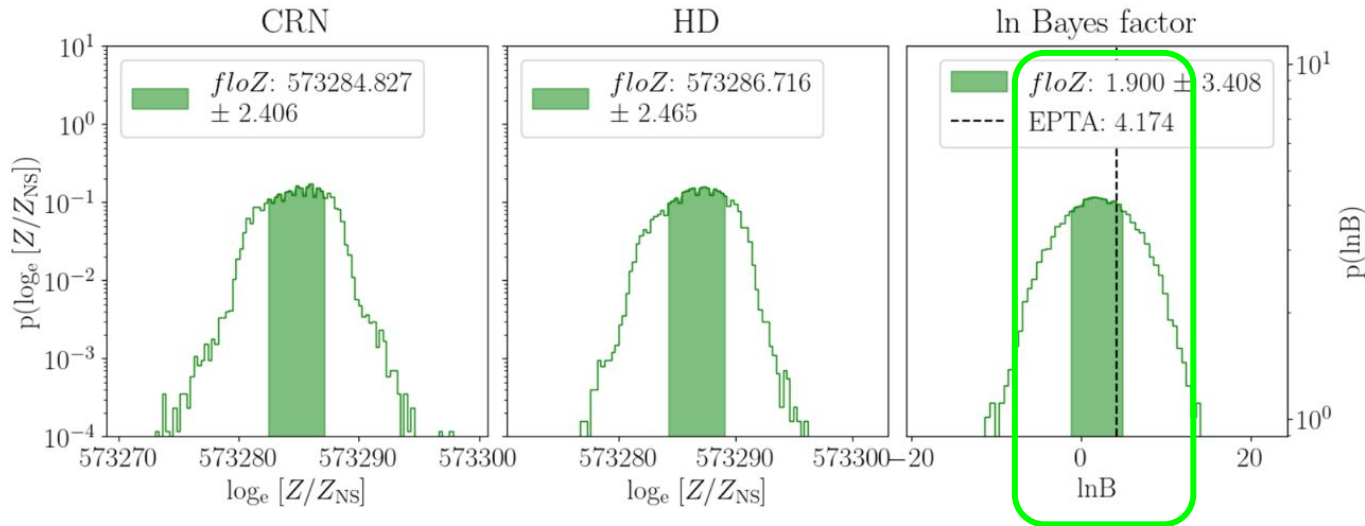
# Applications: Pulsar Timing Array

Bayes factor in favor of the presence of Hellings-Downs relation in EPTA data

70 dimensional samples, with  $1e5$  samples.

$floZ$  estimates is compatible with EPTA results within the  $1\sigma$  uncertainties.

Very non-gaussian distribution  $\rightarrow$  Need more samples (ongoing analysis)



Samples provided by the EPTA collaboration

# Convergence Test

How do we know that the flow is correct?

# Convergence Test

How do we know that the flow is correct?

*$\mu - \sigma$  plot*

Mapping to Latent space should be a gaussian

